

*unmoderated
version*

Machine Learning Course - CS-433

Expectation-Maximization Algorithm

Nov 30, 2022

Martin Jaggi
Last updated on: November 28, 2022

credits to Mohammad Emtiyaz Khan & Rüdiger Urbanke

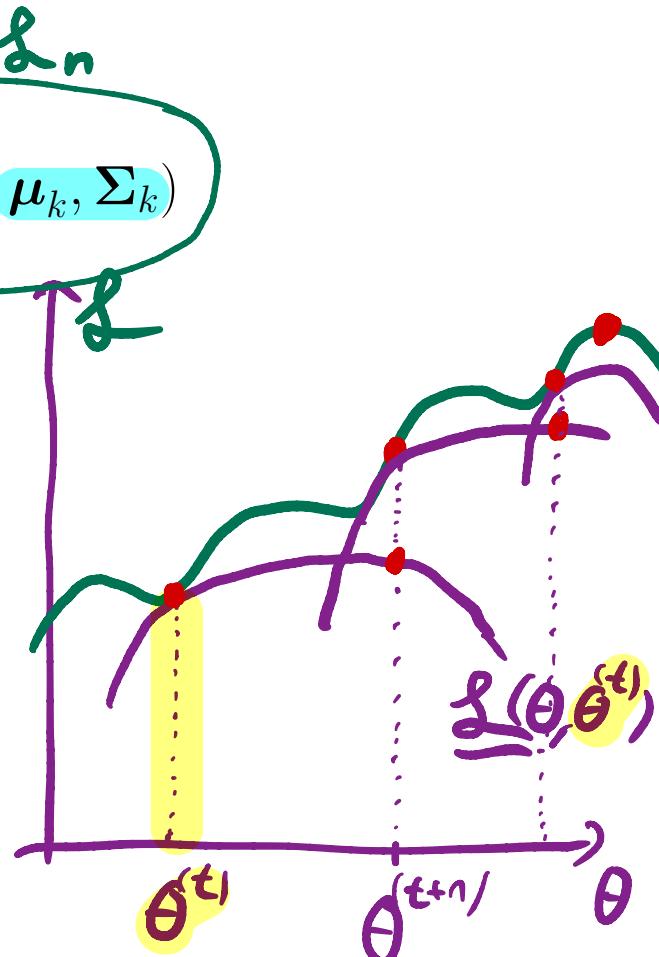


Motivation

Computing maximum likelihood for Gaussian mixture model is difficult due to the log outside the sum.

$$\max_{\theta} \mathcal{L}(\theta) := \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Expectation-Maximization (EM) algorithm provides an elegant and general method to optimize such optimization problems. It uses an iterative two-step procedure where individual steps usually involve problems that are easy to optimize.



EM algorithm: Summary

Start with $\theta^{(1)}$ and iterate:

1. **Expectation step:** Compute a lower bound to the cost such that it is tight at the previous $\theta^{(t)}$:

$$\begin{aligned} \mathcal{L}(\theta) &\geq \underline{\mathcal{L}}(\theta, \theta^{(t)}) \text{ and} \\ \mathcal{L}(\theta^{(t)}) &= \underline{\mathcal{L}}(\theta^{(t)}, \theta^{(t)}). \end{aligned}$$

form $\underline{\mathcal{L}}(\cdot, \theta^{(t)})$

2. **Maximization step:** Update θ :

$$\theta^{(t+1)} = \arg \max_{\theta} \underline{\mathcal{L}}(\theta, \theta^{(t)}).$$

maximize $\underline{\mathcal{L}}(\theta, \theta^{(t)})$

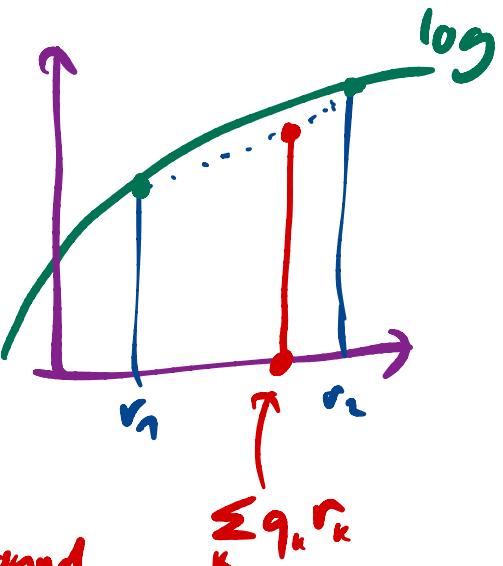
\Leftrightarrow Convexity of $-\log$ Concavity of \log

Given non-negative weights q s.t.

$\sum_k q_k = 1$, the following holds for any $r_k > 0$:

\Leftrightarrow Jensen's inequality

$$\log \left(\sum_{k=1}^K q_k r_k \right) \geq \sum_{k=1}^K q_k \log r_k$$



The expectation step

① $\underline{\mathcal{L}}_n$

$$\log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

lower bound

$$\geq \sum_{k=1}^K q_{kn} \log \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{q_{kn}}$$

② with equality when,

$$q_{kn} = \frac{\pi_k^{(t)} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k^{(t)}, \boldsymbol{\Sigma}_k^{(t)})}{\sum_{k=1}^K \pi_k^{(t)} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k^{(t)}, \boldsymbol{\Sigma}_k^{(t)})}$$

$\underline{\mathcal{L}}_n(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$

This is not a coincidence.

$$\underline{\mathcal{L}}_n(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = \underline{\mathcal{L}}_n(\boldsymbol{\theta}^{(t)})$$

$$\begin{aligned}
 &= \sum_k q_{kn}^{(t)} \log \left(\dots \cdot r_k \dots \right) \\
 &= \sum_{k=1}^K \frac{q_{kn}^{(t)} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k^{(t)}, \boldsymbol{\Sigma}_k^{(t)})}{\sum_{k=1}^K \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k^{(t)}, \boldsymbol{\Sigma}_k^{(t)})} \log \frac{\pi_k^{(t)} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \\
 &= \log \sum_k \pi_k^{(t)} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad 1 \\
 &= \underline{\mathcal{L}}_n(\boldsymbol{\theta}^{(t)})
 \end{aligned}$$

$$\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = e^{-(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) / 2}$$

The maximization step

Maximize the lower bound w.r.t. θ .

$$\underline{\mathcal{L}}_n(\theta, \theta^t)$$

$$\max_{\theta} \sum_{n=1}^N \sum_{k=1}^K q_{kn}^{(t)} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)] \xrightarrow{-\log(q_{kn})} \log \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \dots)}{q_{kn}}$$

Differentiating w.r.t. $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k^{-1}$, we can get the updates for $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$.

$$\boldsymbol{\mu}_k^{(t+1)} := \frac{\sum_n q_{kn}^{(t)} \mathbf{x}_n}{\sum_n q_{kn}^{(t)}}$$

$$\nabla_{\boldsymbol{\mu}_k} \underline{\mathcal{L}}(\theta, \theta^t) \stackrel{!}{=} 0$$

$$\boldsymbol{\Sigma}_k^{(t+1)} := \frac{\sum_n q_{kn}^{(t)} (\mathbf{x}_n - \boldsymbol{\mu}_k^{(t+1)}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{(t+1)})^\top}{\sum_n q_{kn}^{(t)}}$$

$$\nabla_{\boldsymbol{\Sigma}_k} \underline{\mathcal{L}}(\theta, \theta^t) \stackrel{!}{=} 0$$

For π_k , we use the fact that they sum to 1. Therefore, we add a Lagrangian term, differentiate w.r.t. π_k and set to 0, to get the following update:

$$\pi_k^{(t+1)} := \frac{1}{N} \sum_{n=1}^N q_{kn}^{(t)}$$

$$\nabla_{\pi_k} \tilde{\mathcal{L}}(\theta, \theta^t) \stackrel{!}{=} 0$$

constraint
 $\sum \pi = 1$ constraint

$$\tilde{\mathcal{L}} := \underline{\mathcal{L}}_n + \beta (\sum \pi_k - 1)$$

+ k-Means as Summary of EM for GMM a special case

Initialize $\mu^{(1)}, \Sigma^{(1)}, \pi^{(1)}$ and iterate between the E and M step, until $\mathcal{L}(\theta)$ stabilizes.

1. E-step: Compute assignments $q_{kn}^{(t)}$:

$$q_{kn}^{(t)} := \frac{\pi_k^{(t)} \mathcal{N}(\mathbf{x}_n | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{k=1}^K \pi_k^{(t)} \mathcal{N}(\mathbf{x}_n | \mu_k^{(t)}, \Sigma_k^{(t)})}$$

$\propto c \cdot \exp(-\|\mathbf{x}_n - \mu_k\|^2 / \sigma^2)$

\dots

$\begin{cases} 1 & \text{k closest to } \mathbf{x}_n \\ 0 & \text{otherwise} \end{cases}$

2. Compute the marginal likelihood (cost).

$$\mathcal{L}(\theta^{(t)}) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k^{(t)} \mathcal{N}(\mathbf{x}_n | \mu_k^{(t)}, \Sigma_k^{(t)})$$

3. M-step: Update $\mu_k^{(t+1)}, \Sigma_k^{(t+1)}, \pi_k^{(t+1)}$.

$$\mu_k^{(t+1)} := \frac{\sum_n q_{kn}^{(t)} \mathbf{x}_n}{\sum_n q_{kn}^{(t)}}$$

$q_{kn} = z_{kn}$

= k-means update

$$\Sigma_k^{(t+1)} := \frac{\sum_n q_{kn}^{(t)} (\mathbf{x}_n - \mu_k^{(t+1)}) (\mathbf{x}_n - \mu_k^{(t+1)})^\top}{\sum_n q_{kn}^{(t)}}$$

$$\pi_k^{(t+1)} := \frac{1}{N} \sum_n q_{kn}^{(t)}$$

points assigned to cluster

k-Means

If we let the covariance be diagonal i.e. $\Sigma_k := \sigma^2 \mathbf{I}$, then EM algorithm is same as K-means as $\sigma^2 \rightarrow 0$.

and spherical

σ : width of cluster

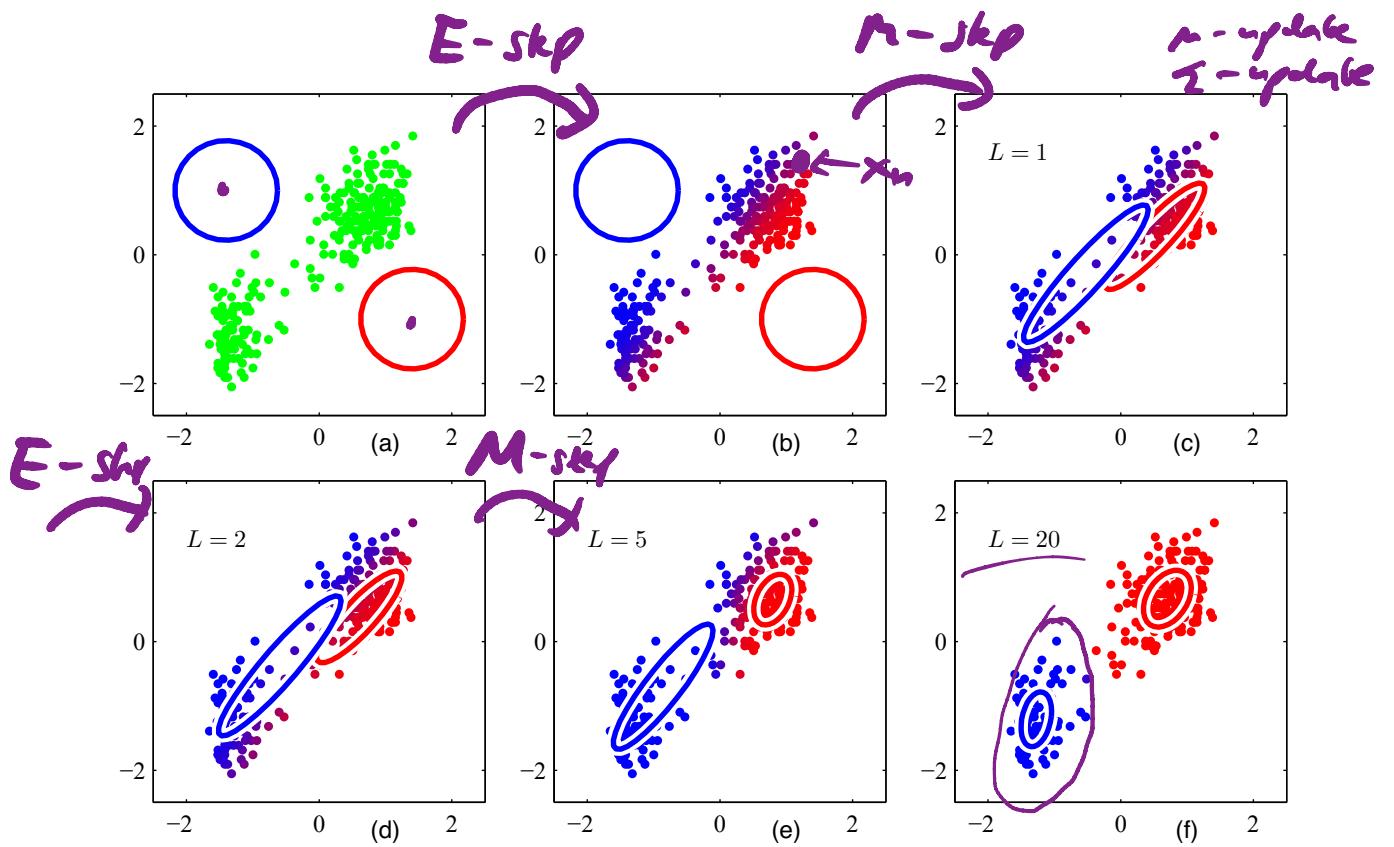


Figure 1: EM algorithm for GMM

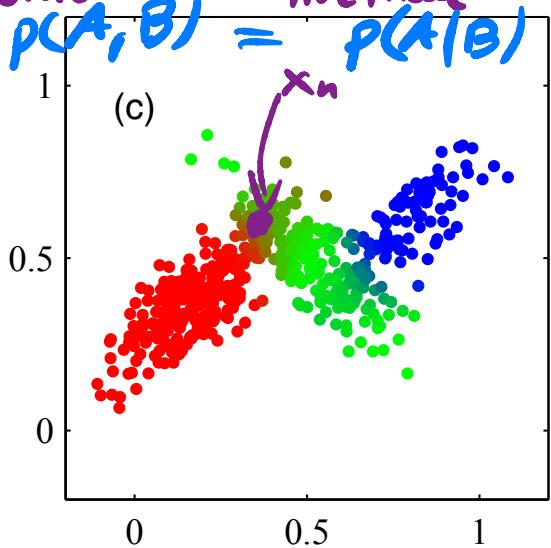
Posterior distribution

We now show that $q_{kn}^{(t)}$ is the posterior distribution of the latent variable, i.e. $q_{kn}^{(t)} = p(z_n = k | \mathbf{x}_n, \theta^{(t)})$

$$p(\mathbf{x}_n, z_n | \theta) = p(\mathbf{x}_n | z_n, \theta) p(z_n | \theta) = p(z_n | \mathbf{x}_n, \theta) p(\mathbf{x}_n | \theta)$$

joint $p(A, B) = p(A|B) \cdot p(B) = p(B|A) \cdot \frac{p(A)}{p(B)}$

Bayes Rule



$$p(z_n = k | \mathbf{x}_n, \theta) = \frac{\text{prior} \cdot \text{likelihood}}{\text{marginal (lik.)}}$$

$$= \frac{\pi_k \cdot N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_k \pi_k \cdot N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}$$

$$= q_{kn}$$

EM in general

Given a general joint distribution $p(\mathbf{x}_n, z_n | \boldsymbol{\theta})$, the marginal likelihood can be lower bounded similarly:

The EM algorithm can be compactly written as follows:

$$\boldsymbol{\theta}^{(t+1)} := \arg \max_{\boldsymbol{\theta}} \sum_{n=1}^N \mathbb{E}_{p(z_n | \mathbf{x}_n, \boldsymbol{\theta}^{(t)})} [\log p(\mathbf{x}_n, z_n | \boldsymbol{\theta})]$$

Expectation over z (the marginal)

maximization

Another interpretation is that part of the data is missing, i.e. (\mathbf{x}_n, z_n) is the “complete” data and z_n is missing. The EM algorithm averages over the “unobserved” part of the data.